

An Introduction to the LDCM Grid Prototype



Jeff Lubelczyk (586)
& Beth Weinstein (586)

January 14, 2005



LDCM Grid Prototype Introduction

- A Grid infrastructure allows scientists at resource-poor sites access to remote resource-rich sites
 - Enable greater scientific research
 - Maximize existing resources
 - Limit expenses
- The objective of the LDCM Grid Prototype (LGP) is to **assess** the applicability and effectiveness of a **data grid** to serve as the **infrastructure** for research scientists to generate virtual Landsat-like data products



LGP Key POCs

■ Sponsors

- LDCM - Bill Ochs, Matt Schwaller
- Code 500/580 - Peter Hughes, Julie Loftis

■ LGP Team members

- Jeff Lubelczyk (Lead)
- Gail McConaughy (SDS Lead Technologist)
- Beth Weinstein (Software Lead)
- Ben Kobler (Hardware, Networks)
- Eunice Eng (Software Dev, Data)
- Valerie Ward (Software Dev, Apps)
- Ananth Rao ([SGT] Software Arch/Dev, Grid Expertise)
- Brooks Davis ([Aerospace Corp] Globus/Grid Admin Expert)
- Glenn Zenker ([QSS] System Admin)

■ USGS

- Stu Doescher (Mgmt)
- Chris Doescher (POC)
- Mike Neiers (Systems Support)

■ Science Input

- Jeff Masek, 923 (Blender)
- Robert Wolfe, 922 (Blender, Data)

■ LDCM Prototype Liaison

- Harper Prior [SAIC]

■ CEOS grid working group (CA)

- Ken McDonald
- Yonsook Enloe [SGT]



High Level Schedule

■ Major Milestones

- 12/03 - Prototype start
- 6/04 - Demo of Capability 1 grid infrastructure
 - Demonstrate simple file transfers and remote application execution at multiple GSFC labs and USGS EDC
 - Ready to build application on top of basic infrastructure
- 12/04 - Demo of Capability 1
 - Provide and demonstrate a grid infrastructure that enables a user program to access and process remote heterogeneous instrument data at multiple GSFC labs and USGS EDC
- 3/05 - Demo of Capability 2 grid infrastructure
 - Demonstrate file transfers and remote application execution at multiple GSFC labs, USGS EDC, and ARC/GSFC commodity resources to assess scalability
- 6/05 - Demo of Capability 2
 - Enable the data fusion (blender) algorithm to obtain datasets, execute, and store the results on any resource within the Virtual Organization (GSFC labs, USGS EDC, ARC/GSFC)



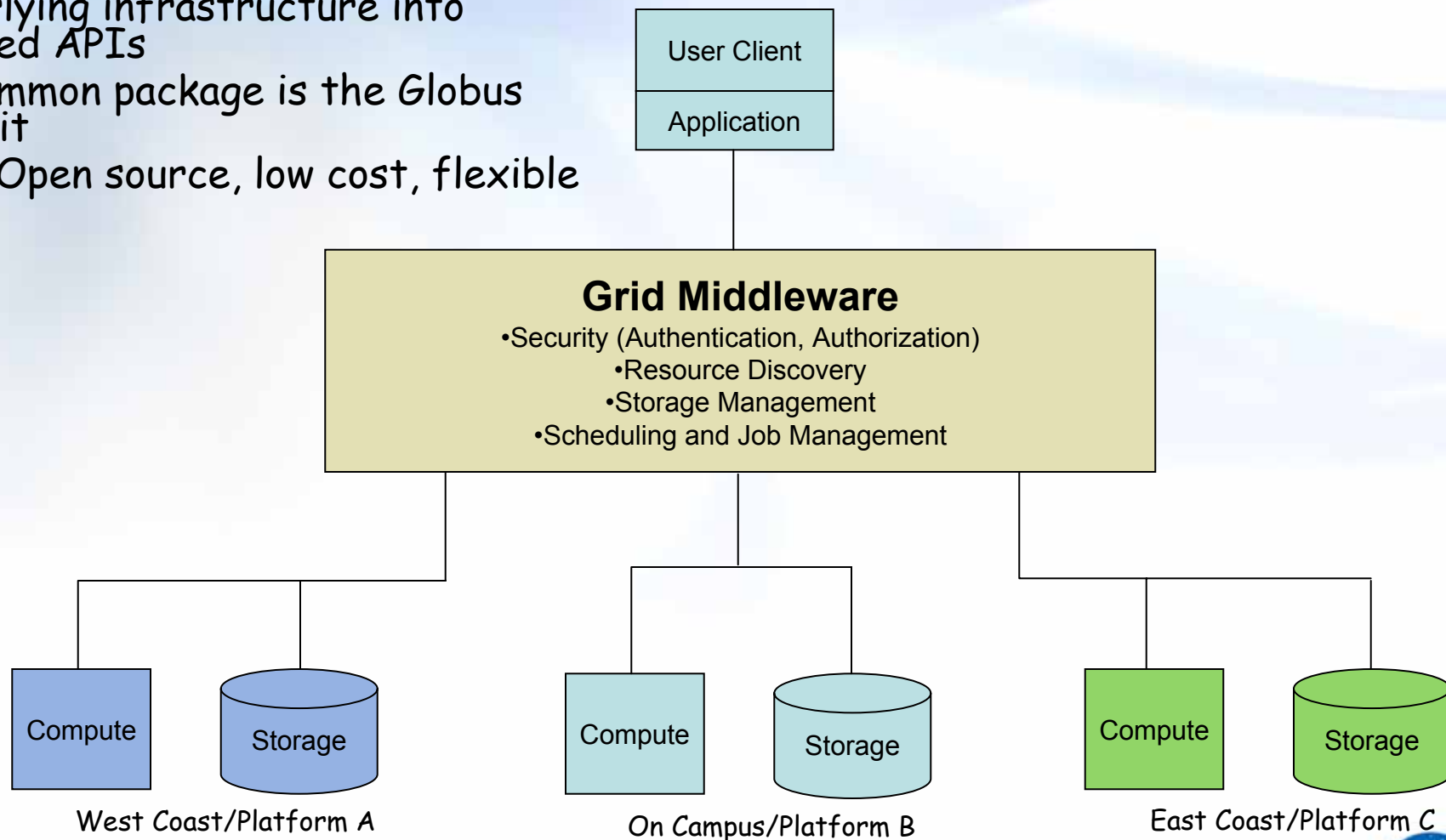
What is a data grid?

- In an article titled "Anatomy of the Grid," Ian Foster of Argonne National Labs suggests the following (2000):
 - "The sharing that we are concerned with is not primarily file exchange but rather **direct access to computers, software, data, and other resources**, as is required by a range of collaborative problem solving and resource-brokering strategies emerging in industry, science, and engineering. This sharing is, necessarily, **highly controlled**, with resource providers and consumers **defining clearly and carefully just what is shared**, **who** is allowed to share, and the **conditions** under which sharing occurs. A set of individuals and/or institutions defined by such sharing rules form what we call a **virtual organization**."
- He further suggests the following criteria:
 - Coordinates resources that are not subject to centralized control
 - Uses standard, open, general purpose protocols and interfaces
 - Otherwise, its an application specific system
 - Delivers nontrivial quality of service
 - Allows resources to be used in a coordinated fashion to deliver varying levels of service



Grid - A Layer of Abstraction

- Grid middleware packages the underlying infrastructure into defined APIs
- A common package is the Globus Toolkit
 - Open source, low cost, flexible





What the current LGP data grid utilizes

■ Security Infrastructure

- Globus Gate Keeper
 - Authentication (PKI)
 - Authorization

■ Resource Discovery

- Monitoring and Discovery Service (MDS) [LDAP like]

■ Storage Management and Brokering

- Metadata catalogs
- Replica Location Service
 - Allows use of logical file names
 - Physical locations are hidden
- Storage Resource Management
 - GridFTP
 - Retrieves data using physical file names
 - Data formats and subsetting

■ Job Scheduling and Resource Allocation

- GRAM (Globus Resource Allocation Manager) -- Provides a single common API for requesting and using remote system resources

Globus Toolkit 2.4.3

Globus
Gate
keeper

GridFTP

GRAM

Note: Portions of the
Globus Toolkit used in
Capability 1





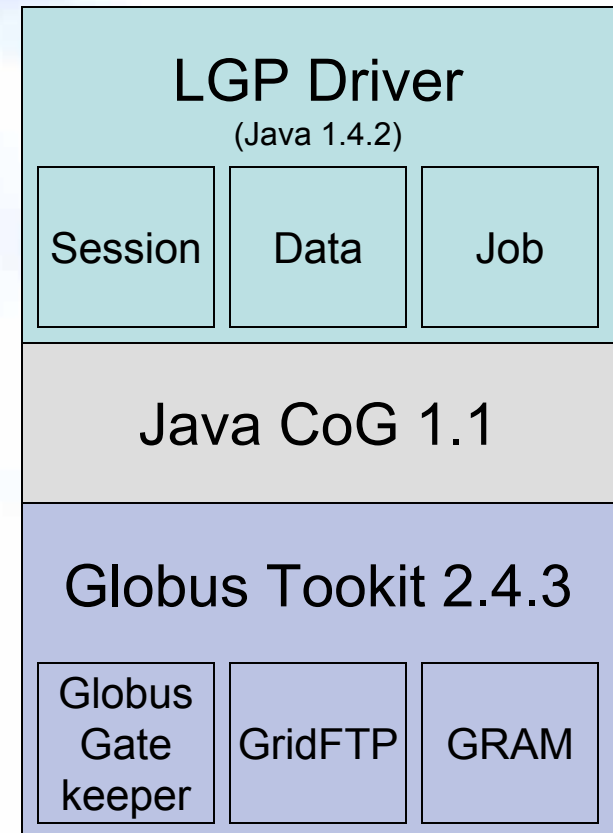
Capability 1 Software Framework

■ LDCM Grid Prototype (LGP) Driver

- Provides a generic software system architecture based on Globus services
- LGP Driver high-level services
 - **Session Manager** - grid session initiation and user authentication using proxy certificates
 - **Data Manager** - file transfer using GridFTP
 - **Job Manager** - job submission and status in a grid environment

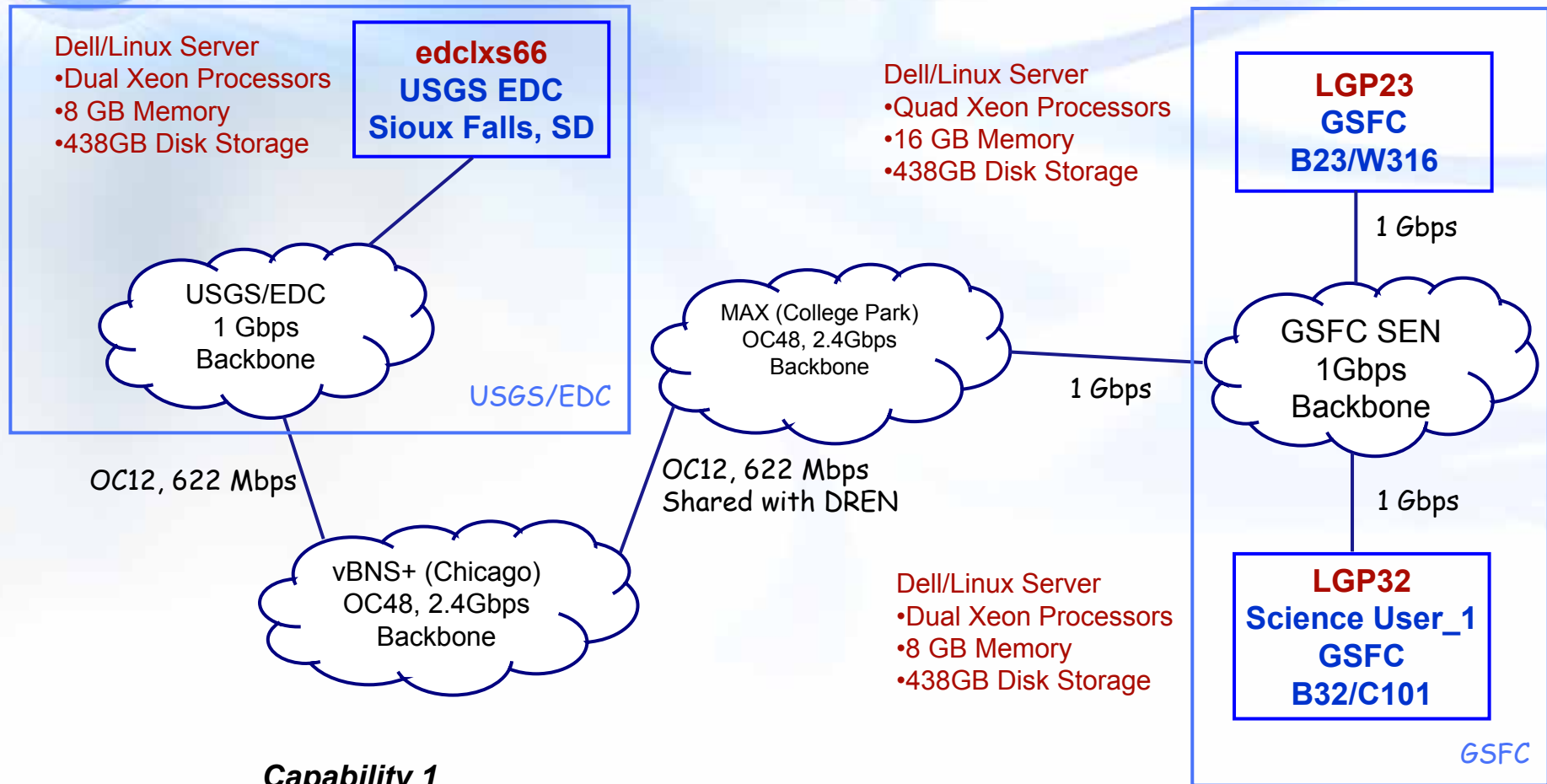
■ Utilizes the Java Commodity Grid Kits (CoGs)

- Supplies a layer of abstraction from underlying Globus services
- Simplifies the programming interface





Capability 1 Virtual Organization



SEN: Science and Engineering Network
MAX: Mid-Atlantic Crossroads
DREN: Defense Research and Engineering Network
vBNS+: Very high Performance Network Service





4 Scenarios to Illustrate Grid Flexibility

- **Data Services** (*Move application to data*)
 - Transfer MRT components to remote hosts and process the data remotely, sending the results back to the science facility
- **Batch Execution** (*Parallel computing*)
 - Demonstrate the execution of MRT components in a parallel batch environment
- **Local Processing** (*User prefers to process locally*)
 - Transfer the selected data sets to the science user site for processing
- **Third Party Processing** (*No local resource usage*)
 - Perform a third party data transfer and process the data remotely

Grid flexibility maximizes science resources





Next Steps -- Capability 2

■ Capability 2 (C2)

- Integrate with the Blender team
 - Collaborate to identify meaningful C2 data sets
 - Demonstrate blender algorithm
- Assess Grid performance
 - Expand the VO to include ARC supercomputing if available
 - Performance Goals
 - Demonstrate the processing of 1 day's worth of data in the grid environment (~250 scenes)
- Grid Workflow -- increase automation
 - Our current capabilities allow us to submit jobs only to a specified resource
 - The goal of the next phase will be to provide the ability to submit a job to the "Grid" Virtual Organization
 - Grid resource management
 - Scheduling policy
 - Reliable job completion
 - Checkpointing and job migration
 - Leverage wasted cpu cycles

